

アミノ酸配列からタンパク質の機能を探る ファミリー、ドメイン、モチーフ

清水謙多郎

shimizu@bi.a.u-tokyo.ac.jp

実習の資料について

- 「バイオインフォマティクス 実習」で検索
- <http://lecture.ecc.u-tokyo.ac.jp/~ashimizu/>

バイオインフォマティクス実習 (清水謙多郎)

ID: bioinfo
パスワード: 5455

ログインした状態で、講義資料からのデータのダウンロードがスムーズです。

モチーフ 2

タンパク質の機能予測

分子機能の推定

- アミノ酸配列
- ホモロジー検索
- ドメイン
- モチーフ
- ドメイン検索、ファミリー検索
- モチーフ検索

• 配列 (特徴の) 類似性

• 構造 (特徴の) 類似性

• 進化に関する情報

• ゲノムコンテキスト

• 遺伝子発現の関連性

• 代謝経路上の位置

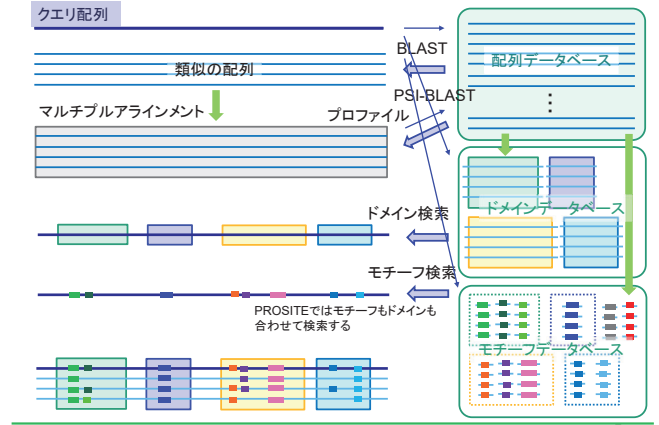
...

生物学的プロセスの推定

- 翻訳後修飾
- 細胞内局在
- 翻訳後修飾予測
- 細胞内局在予測、膜貫通予測
- 他の分子との相互作用
- 相互作用予測、相互作用部位予測

モチーフ 3

タンパク質の配列からの機能予測



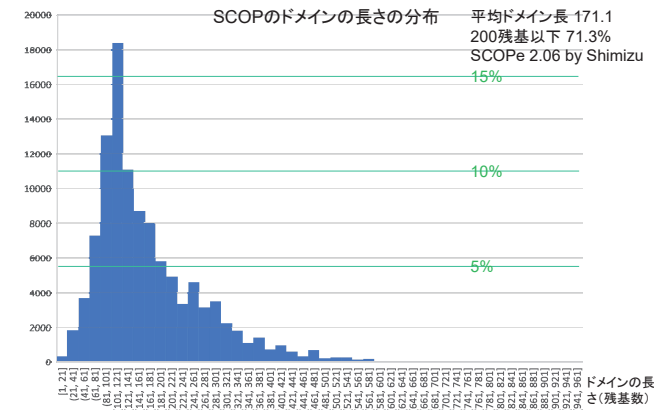
5

タンパク質のドメイン

- ドメイン (domain): タンパク質構造を構成する構造的あるいは機能的にまとまった単位
 - タンパク質は1つまたは複数のドメインから構成される
 - 多くは200残基程度からそれ以下
 - G. A. Petsko, D. Ringe. *Protein Structure and Function*, New Science Press, 2004.
- 各ドメインには、共通した配列・構造特徴および機能をもつものがあり、それらは「同じドメイン」とみなされる
- 「同じドメイン」が複数のタンパク質に現れることが多い
- 複数のドメインから構成されるマルチドメインタンパク質では、各ドメインに対応する機能を合わせもつことが多い → ドメインの構成によって、タンパク質の機能が決まる

モチーフ 6

参考 タンパク質のドメイン



GO 7

BLAST検索の結果

ドメインの情報が表示される

Putative conserved domains have been detected, click on the image below for detailed results.

Distribution of the top 100 Blast Hits on 100 subject sequences

参考 タンパク質のドメイン

- 真核生物のタンパク質の65%、原核生物のタンパク質の40%がマルチドメイン構成である
 - Ekman D, et al. "Multi-domain proteins in the three kingdoms of life: orphan domains and other unassigned regions". *Journal of Molecular Biology*. 348: 231-243 (2005). doi:10.1016/j.jmb.2005.02.007.
- ドメインは、それ自身、独立してフォールドできるような構造で、真核生物のマルチドメインタンパク質の1つのドメインが、原核生物で単独のタンパク質で存在する例が見られる
 - Davidson JN, Chen KC, Jamison RS, Musmanno LA, Kern CB (March 1993). "The evolutionary history of the first three enzymes in pyrimidine biosynthesis". *BioEssays*. 15: 157-164. doi:10.1002/bies.950150303.

三機能プリン合成タンパク質アデニン-3
GAR合成酵素 (GARS)、AIR合成酵素 (AIRS)、GARホルミル基転移酵素 (GART)



モチーフ 8

ドメインの例

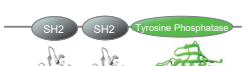
プロテインチロシンキナーゼ Src



UniProtKB P12931
PDB 2h8h

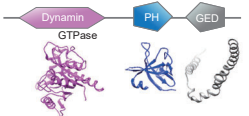
SH2: Src homology 2 domain
SH3: Src homology 3 domain

プロテインチロシンホスファターゼ



UniProtKB Q06124
PDB 2shp

ダイナミン



UniProtKB Q9UQ16
PDB 5a3f

PH: Pleckstrin homology domain
GED: Dynamin GTPase effector domain

モチーフ 9

Rieske型鉄硫黄クラスター結合部位

Rieske型鉄硫黄クラスター結合部位の配列アラインメント



PROSITE PS51296

Rieske [2Fe-2S] 鉄硫黄ドメインのプロファイル

ドメイン全体にわたる配列特徴

C-x-H-x(15-17)-C-x-x-H

鉄結合部位のモチーフ

モチーフの局所的な配列特徴

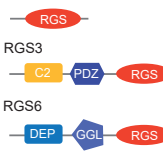
モチーフ 13

ファミリーとドメイン

タンパク質ファミリー

- 進化的類縁関係をもつタンパク質のグループ
 - 共通の機能をもつタンパク質のグループ
- 「ドメインの構成によって、タンパク質の機能が決まる」
- ファミリーはドメインの構成によって決まる
 - 同じドメインが異なるファミリーに現れることがある

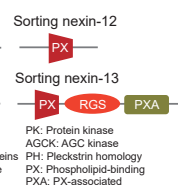
Gタンパク質シグナル伝達調節因子 (RGS)



βアドレナリン受容体キナーゼ (BARK)



ソーティングネキシン



RGS: Regulator of G-protein signaling
C2: Membrane targeting
PDZ, DEP: Involved in the signaling proteins
GGL: Found in the gamma subunit of the heterotrimeric G protein complex, etc

PK: Protein kinase
AGCK: AGC kinase
PH: Pleckstrin homology
PX: Phospholipid-binding
PXA: PX-associated

モチーフ 10

ファミリーとドメイン

- タンパク質の機能を探るのに、構成するドメインを同定することは重要
- ドメインとその機能、タンパク質が構成するドメインをデータベースに登録
 - 配列の類似性を手がかりとして、ドメインを同定し、機能を推定する
 - ドメインを同定することは、タンパク質のファミリーを同定することにつながる

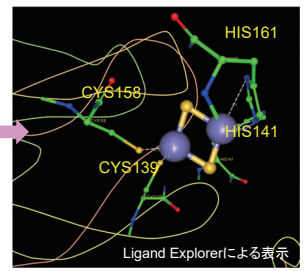


ドメイン、ファミリーのデータベースと高感度の検索

モチーフ 11

鉄硫黄クラスター周辺の構造

Cytochrome bc1 (UniProtKB: UCRI_BOVIN, PDB: 1L0L Eチェーン)
鉄硫黄クラスター (2Fe-2S) 周辺の構造



結合部位 (結合サイト)

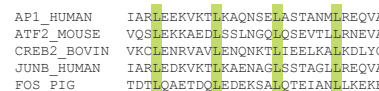
モチーフ 14

Leucine Zipper

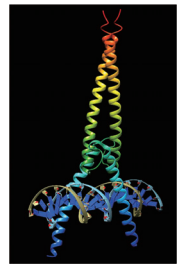
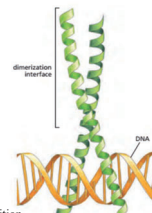
DNAに結合する転写因子に見られる、αヘリックスの2量体

L-x(6)-L-x(6)-L-x(6)-L

PROSITE PS00029



N末端側に塩基性アミノ酸が見られる領域



OmoMYC bound to double-stranded DNA
PDBID: 5i50

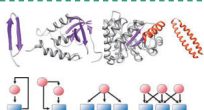
Molecular Biology of the Cell, 6 edition

モチーフ 15

モチーフ

- 複数の塩基配列やアミノ酸配列に共通に見られる、保存された(短い)配列パターン
 - 特定の機能、性質に関係する
 - 塩基配列 → 転写因子結合部位、TATAボックス、スプライシング部位(GT-AGモチーフ)など
- アミノ酸配列のモチーフ
 - 酵素の活性部位、他の分子との相互作用部位、翻訳後修飾、ドメインを特徴づける配列パターンなどに関係
 - (配列全体の類似度が低くても) 特定の機能、性質に関係する部分は強く保存されている傾向にある

タンパク質の立体構造 → 構造モチーフ



ネットワーク → ネットワークモチーフ



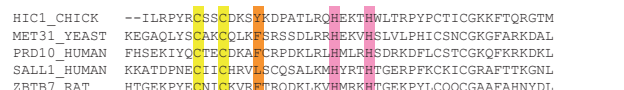
モチーフ 12

Zinc Finger

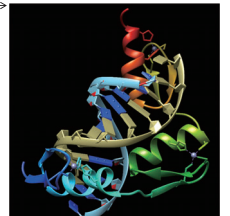
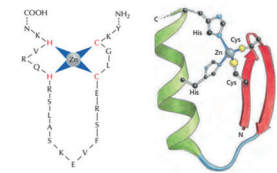
DNAに結合する転写因子に見られる

C-x(2,4)-C-x(3)-[LIVMFYWCI]-x(8)-H-x(3,5)-H

PROSITE パターン PS00028
プロファイル PS0157



シート シート ヘリックス



Molecular Biology of the Cell, 6 edition

モチーフ 16

Zif268-DNA complex, PDBID: 1zaa

PROSITEデータベース

- タンパク質のファミリー、ドメイン、機能部位のデータベース
- それらを特徴づける配列パターン(モチーフ)を登録し、その検索機能をもつ
- ファミリーに特徴的な機能やドメインの構造を記述したドキュメンテーションエントリーと、モチーフの表現を登録したエントリーをもつ
- EBIのExPaSy Proteomics Serverのデータベースの一つ
- <https://prosite.expasy.org/>

PROSITEデータベースの利用

参考

PROSITEの検索(3)

- PDOC00611 XPA protein signatures
- PDOC01915 YDQ domain profile
- PDOC01972 Yippee domain profile
- PDOC00683 Y6C-like domain profile
- PDOC01912 ZAD domain profile
- PDOC01936 Zinc finger A20-type profile
- PDOC01939 Zinc finger AN1-type profile
- PDOC02919 Zinc finger B-box-type profile
- PDOC02908 Zinc finger BED-type profile
- PDOC01913 Zinc finger Bix-type profile
- PDOC01799 Zinc finger C2H2 AKAP95-type profile
- PDOC00028 Zinc finger C2H2-type domain signature and profile**
- PDOC01944 Zinc finger C2H2 LNRQ-type profile
- PDOC01803 Zinc finger C2H2 RNF-type profile
- PDOC01807 Zinc finger C2H2 Baccalarius (BV)-type profile
- PDOC01903 Zinc finger C3H1-type profile
- PDOC01896 Zinc finger C4H2-type profile
- PDOC01910 Zinc finger C4H2-FHQ-type profile
- PDOC01811 Zinc finger C4H2-HVEP-type profile
- PDOC01801 Zinc finger C4H2-NQA-type profile
- PDOC01918 Zinc finger C4H2-type profile
- PDOC01802 Zinc finger C4H2C-type profile
- PDOC01802 Zinc finger CHC1-U11-48k-type profile
- PDOC01926 Zinc finger CHY-type and C1CHY-type profile
- PDOC01959 Zinc finger CW-type profile
- PDOC01959 Zinc finger CXC-type profile
- PDOC01901 Zinc finger DNL-type profile
- PDOC01902 Zinc finger DMA-type profile
- PDOC00884 Zinc finger Dot-type signature and profile
- PDOC01924 Zinc finger FCS-type profile
- PDOC01919 Zinc finger FLZ-type profile
- PDOC01918 Zinc finger FYVE-FYVE-related type profile
- PDOC01803 Zinc finger HT-type profile
- PDOC00895 Zinc finger MYND-type signature and profile
- PDOC00616 Zinc finger PHD-type signature and profile
- PDOC01902 Zinc finger RAG1-type profile
- PDOC01929 Zinc finger RING-CH1-type profile
- PDOC00449 Zinc finger RING-type signature and profile
- PDOC01919 Zinc finger RIT-type profile
- PDOC01909 Zinc finger RanBP2-type signature and profile
- PDOC01914 Zinc finger SSP-type profile
- PDOC01811 Zinc finger SNAH-type profile
- PDOC01444 Zinc finger SP-DNA-Chrom-remode

PDOC00028 「zinc finger C2H2-type domain signature and profile」を選択

参考

PROSITEの検索(4)

参考

PROSITEの検索(1)

参考

PROSITEの検索(5)

参考

PROSITEの検索(2)

参考

PROSITEの検索(6)

PROSITEの検索(12)

SARS CoV2とSARS CoVの比較(2)

```

CoV2_Mpro SGFRKMAFFSGRVEGCMVQVTCGTTLNLGLDLDVYVCPFRVICTSEMLNPNFYEDLLIRKSNHFLVQAGNVQLRVIGH
SARCoV2_Mpro SGFRKMAFFSGRVEGCMVQVTCGTTLNLGLDLDVYVCPFRVICTSEMLNPNFYEDLLIRKSNHFLVQAGNVQLRVIGH
CoV2_Mpro SMQNCVLRKLVDTANPKTKPKYKFRVIRIQGQTFSLVACYNGSPGVYQCAMRPNFTIKGSFLNGS;GSGVFNIDYDCVSEFC
SARCoV2_Mpro SMQNCVLRKLVDTANPKTKPKYKFRVIRIQGQTFSLVACYNGSPGVYQCAMRPNFTIKGSFLNGS;GSGVFNIDYDCVSEFC
CoV2_Mpro YMHMELPTGVHAGTDLGKGFYGFVDRQTAQAAGTDTTITLVNLMVLAAYVINGDRWFLNRFITLNDFNLVAMKYNVE
SARCoV2_Mpro YMHMELPTGVHAGTDLGKGFYGFVDRQTAQAAGTDTTITLVNLMVLAAYVINGDRWFLNRFITLNDFNLVAMKYNVE
CoV2_Mpro PLTQDQVVDILGPLSAQTGIAVLDMCA;SLKELLQNGMGRITLGLSALLEDEFTFFDVRVRCQSGVTFQ
SARCoV2_Mpro PLTQDQVVDILGPLSAQTGIAVLDMCA;SLKELLQNGMGRITLGLSALLEDEFTFFDVRVRCQSGVTFQ
    
```

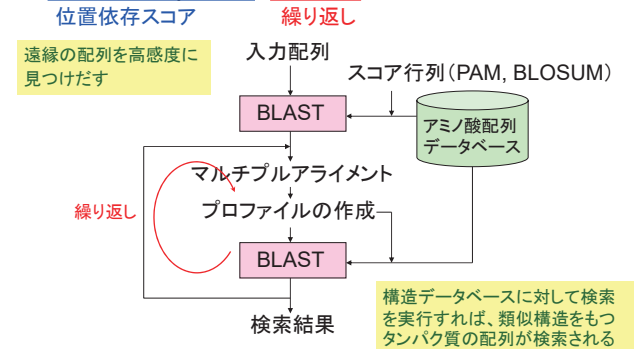
GO 37

PROSITEの検索(13)

モチーフ 34

PSI-BLAST

Position-Specific Iterated BLAST



モチーフ 38

PROSITEの検索(14)

GO 35

参考 PSI-BLASTのプロファイル

- 配列のマルチプルアライメントの各位置(カラム) i において、
 - $n(i, j)$: その位置のアミノ酸 j の出現度数
 - その位置のアミノ酸 j の出現頻度 $f(i, j)$ は、

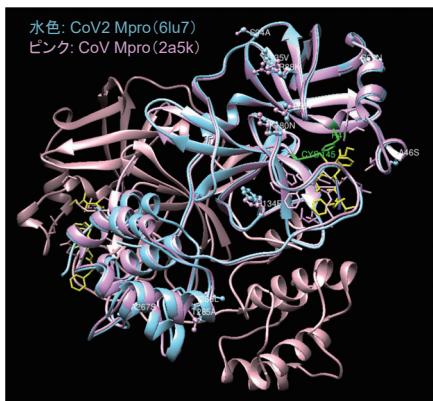
$$f(i, j) = \frac{n(i, j)}{\sum_j n(i, j)}$$
- PSI-BLASTのプロファイル $S(i, j)$ の計算
 - $P(j)$ はバックグラウンドでのアミノ酸 j の出現頻度
 - $Q(i, j)$ は、その位置に現れると見られるアミノ酸 j の出現頻度

$$Q(i, j) = \frac{\alpha f(i, j) + \beta g(i, j)}{\alpha + \beta}$$
 - α は、そのカラムに出現するギャップも含めたアミノ酸の種類数 - 1、 $\beta = 10$
 - $q(i, j)$ をアミノ酸の置換スコアとすると、

$$g(i, j) = \sum_k \frac{f(i, k)}{P(k)} q(k, j)$$

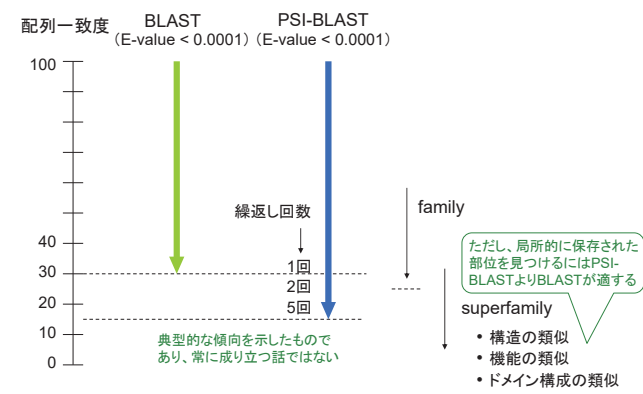
PSI-BLAST 39

SARS CoV2とSARS CoVの比較(1)



GO 36

PSI-BLASTによる検索の基準



モチーフ 40

参考 PSI-BLASTの利用(1)

モチーフ 41

参考 PSI-BLASTの利用(5)

モチーフ 45

参考 PSI-BLASTの利用(2)

モチーフ 42

参考 PSI-BLASTの利用(6)

モチーフ 46

参考 PSI-BLASTの利用(3)

モチーフ 43

参考 PSI-BLASTの利用(7)

モチーフ 47

参考 PSI-BLASTの利用(4)

モチーフ 44

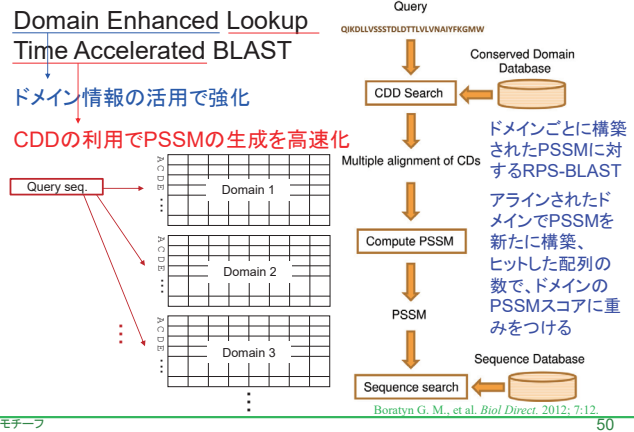
参考 DELTA-BLASTの利用(1)

モチーフ 48

参考 DELTA-BLASTの利用(2)

Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per Ident	Acc Len	Accession	Select for PSI-BLAST iteration 1	Used for newly added PSSM
Acetaminophen Aden complex with Ser domain D (Acetaminophen castellanii)	Acetaminophen	586	586	99%	0.0	92.45%	375	HEFL_A	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Cardiac thin filament decorated with COC1 fragment of cardiac myosin binding protein C (Homo sapiens)	Homo sapiens	579	579	99%	0.0	93.80%	375	SCDA_A	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Cristal structure of Dichoelasma Dichocentrum Aden Complexed With Ca ²⁺ ATP And Human...	Dichoelasma	578	578	99%	0.0	91.13%	375	SLN1_A	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
Cryo-EM structure of a human retinoplasmic adenosine complex at near-atomic resolution (Homo sapiens)	Homo sapiens	578	578	99%	0.0	93.80%	374	SLN1_A	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>

参考 DELTA-BLAST



参考 DALI

- DALI: タンパク質の構造比較の手法およびデータベース
<http://ekhidna2.biocenter.helsinki.fi/dali/>

参考 DALIの検索結果の例(1)

参考 DALIの検索結果の例(2)

参考 DALIの検索結果の例(3)

参考 DALIの結果の例(2)

```
DSP HLLLLLHRRHH-LLEEEELARRHLLLHRRHHHRRHLL11111LEELLHRRHH
Query INKCDIDRDLY-ANNVWSGTTMYPGIADRQKELTlqpsmtKIKIAPPKRYKVV 340
ident | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
Sbjct LKTPFPELVSDILeRGIPLTQGGSLIKGLDFLLKQETG-----LIVISEKPTLAVV 313
DSP HLLLLLHRRHHHLLLEEEELARRHLLLHRRHHHRRHLL11111LEELLHRRHH
```

```
DSP HRRHHRRHLL-LHHLLEHHHHHHH11hhhh1
Query IGGSLASLST-PQQMMWIKTqydeagsgvsvhr 372
ident | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
Sbjct KGAGXVLDFWIKIKKLGAG----- 333
DSP HRRHHRRHLLHLLLEEL
```

Pfam

- Protein families database of alignments and HMMs
- タンパク質のドメインとファミリー、それらの特徴づける配列パターンを表す隠れマルコフモデル(HMM)を登録
- ドメインの立体構造も表示(そのドメインをもつタンパク質の構造)
- 人手によるアノテーション
- <https://pfam.xfam.org/>

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam 34.0 (March 2021, 19179 entries)

The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs). More...

QUICK LINKS YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...
SEQUENCE SEARCH Analyze your protein sequence for Pfam matches
VIEW A PFAM ENTRY View Pfam annotation and alignments
VIEW A CLAN See groups of related entries
VIEW A STRUCTURE Look at the domain organisation of a protein sequence
VIEW A STRUCTURE Find the domains on a PDB structure
KEYWORD SEARCH Query Pfam by keywords
JUMP TO Enter any accession or ID to jump to the page for a Pfam entry or clan, UniProt sequence, PDB structure, etc. Or view the [help page](#) for more information.

Recent Pfam blog posts
Google Research Team blog Deep Learning to Pfam (Posted 24 March 2022)
 We are delighted to announce the first fruits of a collaboration between the Pfam team and a Google Research team led by Dr Lucy Colwell, with Maxime Bessard and David Berger. In 2019, Colwell's team published a preprint describing a new deep learning method that was trained on Pfam data, and which improves upon the [...] [more](#)

Pfam 34.0 is released! (Posted 24 March 2022)
 Pfam 34.0 contains a total of 18,379 families and 645 clans. Since the last release, we have built 935 new families, added 15 families and created 11 new clans. UniProt Reference Proteomes has increased by 21% since Pfam 33.1, and now contains 47 million sequences. Of the sequences that are in reference proteomes, 74.5% have [...] [more](#)

Folding the Protein Universe (Posted 3 March 2022)
 Today signifies the realization of a long-held dream to have the structure of every (well nearly every) family in Pfam. The Pfam and InterPro databases have made available structural models of 6,270 protein families created by Ivan Anishchanka from David Baker's group at the University of Washington in Seattle. The models are made using their [...] [more](#)

Pfamの検索(1)

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam 34.0 (March 2021, 19179 entries)

The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs). More...

QUICK LINKS YOU CAN FIND DATA IN PFAM IN VARIOUS WAYS...
SEQUENCE SEARCH Analyze your protein sequence for Pfam matches
VIEW A PFAM ENTRY View Pfam annotation and alignments
VIEW A CLAN See groups of related entries
VIEW A STRUCTURE Look at the domain organisation of a protein sequence
VIEW A STRUCTURE Find the domains on a PDB structure
KEYWORD SEARCH Query Pfam by keywords
JUMP TO Enter any accession or ID to jump to the page for a Pfam entry or clan, UniProt sequence, PDB structure, etc. Or view the [help page](#) for more information.

Recent Pfam blog posts
Google Research Team blog Deep Learning to Pfam (Posted 24 March 2022)
 We are delighted to announce the first fruits of a collaboration between the Pfam team and a Google Research team led by Dr Lucy Colwell, with Maxime Bessard and David Berger. In 2019, Colwell's team published a preprint describing a new deep learning method that was trained on Pfam data, and which improves upon the [...] [more](#)

Pfam 34.0 is released! (Posted 24 March 2022)
 Pfam 34.0 contains a total of 18,379 families and 645 clans. Since the last release, we have built 935 new families, added 15 families and created 11 new clans. UniProt Reference Proteomes has increased by 21% since Pfam 33.1, and now contains 47 million sequences. Of the sequences that are in reference proteomes, 74.5% have [...] [more](#)

Folding the Protein Universe (Posted 3 March 2022)
 Today signifies the realization of a long-held dream to have the structure of every (well nearly every) family in Pfam. The Pfam and InterPro databases have made available structural models of 6,270 protein families created by Ivan Anishchanka from David Baker's group at the University of Washington in Seattle. The models are made using their [...] [more](#)

Pfamの検索(2)

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Pfam 34.0 (March 2021, 19179 entries)

The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs). More...

「SEQUENCE SEARCH」を選択

配列を入力するための枠が開く

「SEQUENCE SEARCH」を選択
SEQUENCE SEARCH Paste your protein sequence here to find matches with entries.
VIEW A PFAM ENTRY View Pfam annotation and alignments
VIEW A CLAN See groups of related entries
VIEW A STRUCTURE Look at the domain organisation of a protein sequence
VIEW A STRUCTURE Find the domains on a PDB structure
KEYWORD SEARCH Query Pfam by keywords
JUMP TO Enter any accession or ID to jump to the page for a Pfam entry or clan, UniProt sequence, PDB structure, etc. Or view the [help page](#) for more information.

「PKinase」に属する38個のメンバー

「PKinase」に属する38個のメンバー

「PKinase」に属する38個のメンバー

Pfamの検索(3)

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Sequence search results

Show the detailed description of this result page.
 We found 4 Pfam-A matches to your search sequence (3 significant and 1 insignificant)

Show the search options and sequences that you submitted.
[Return to the search form to look for Pfam domains on a new sequence.](#)

Significant Pfam-A Matches
 Data of 20/20 alignments

Family	Description	Entry type	Clan	Envelope	Alignment	HMM	ES	ES score	E-value	Predicted active sites	Show/hide alignment		
SH3_1	SH3 domain	Domain	CL0010	270	270	270	2	258	209.2	1.4e-92	View		
SH3_2	SH3 domain	Domain	CL0010	151	233	151	233	1	77	95.3	1.8e-27	View	
SH3_3	SH3 domain	Domain	CL0010	90	137	90	137	1	48	48	60.2	1.9e-16	View

検出されたドメイン・ファミリーの情報

ドメイン・ファミリーの名前を指定すると、それらの説明が表示される

ドメインのコンセンサス配列とクエリ配列のアラインメントが表示される

Pfamの検索(4)

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Family: SH3_1 (PF001018)

Summary: SH3 domain

Domain organisation
 Class
 Alignments
 HMM logo
 Trees
 Clusters & model
 Species
 Interactions
 Structures
 Jump to...

InterProのエントリー

旧来のPfamのページ

構造の例

外部データベースのリンク

Pfamの検索(5)

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Family: Pkinase_Tyr (PF07714)

Summary: Protein tyrosine kinase

Domain organisation
 Class
 Alignments
 HMM logo
 Trees
 Clusters & model
 Species
 Interactions
 Structures
 Jump to...

Protein tyrosine kinase

構造の例

Pfamの検索(6)

EMBL-EBI HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

Family: Kinase_Tyr (PF07714)

Summary: Kinase_Tyr

Domain organisation
 Class
 Alignments
 HMM logo
 Trees
 Clusters & model
 Species
 Interactions
 Structures
 Jump to...

Clan「PKinase」(protein kinase domain)

「PKinase」に属する38個のメンバー

CD-Searchの結果(1)

The screenshot shows the NCBI Conserved Domains search interface. The query sequence is 'seq1.txt (チロシンキナーゼSRC)'. The search results show several domain hits, with 'SH2_Src' and 'PTKc_11ke' being the most prominent. A red circle highlights the 'SH2_Src' domain hit. The 'List of domain hits' table is as follows:

Name	Accession	Description	Start	End
PTKc_11ke	cd05071	Catalytic domain of the Protein Tyrosine Kinase, Src. PTKs catalyze the transfer of the gamma phosphate of ATP to cytosolic (or non-receptor) PTKs, containing an SH2 domain, and is negatively regulated by phosphorylation at the C-terminal domain, forming of the oncogenic oncoprotein Src. Src is a member of the Src family of the oncogenic protein tyrosine kinases. The SH2 domain is a conserved domain of the Src family of the oncogenic protein tyrosine kinases. The SH2 domain is a conserved domain of the Src family of the oncogenic protein tyrosine kinases.	250-258	504-510
SH2_Src	cd12008	SH2 domain of the Protein Tyrosine Kinase, Src. SH2 is a conserved domain of the Src family of the oncogenic protein tyrosine kinases. The SH2 domain is a conserved domain of the Src family of the oncogenic protein tyrosine kinases.	147-247	306-312
SH2_Src	cd12008	SH2 domain of the Protein Tyrosine Kinase, Src. SH2 is a conserved domain of the Src family of the oncogenic protein tyrosine kinases. The SH2 domain is a conserved domain of the Src family of the oncogenic protein tyrosine kinases.	80-143	400-433

タンパク質の機能情報のレベル

スーパーファミリー	進化的に類縁関係にあるタンパク質のグループ。進化的類縁関係という尺度で最大限に感度を上げてまとめたもの。
クラン (clan)	類似の機能をもつタンパク質のファミリーのグループ。ファミリーをより大きな枠組みで括るのに用いられ、進化は必ずしも関係ない。
ファミリー	進化的な類縁関係が明確で、配列類似度が高いタンパク質のグループ。
サブファミリー	ファミリーを進化的類縁関係、機能により、さらに細かく分類したグループ。ファミリー内の複数の機能を区別して論じるのに用いられる。
ドメイン	タンパク質内の構造的あるいは機能的に独立した部分構造。複数のタンパク質に繰り返し現れる。配列の類似性が見られる。
モチーフ	複数のアミノ酸配列で共通に見られる短い配列の特徴的なパターン。
サイト	高度に保存された短い配列の部分。
シグニチャー (signature)	関連する複数のアミノ酸配列に共通に現れる配列パターンの総称。

CD-Searchの結果(2)

The screenshot shows the NCBI Conserved Protein Domain Family search results for 'PTKc_Src'. The search results show several domain hits, with 'PTKc_Src' being the most prominent. The 'List of domain hits' table is as follows:

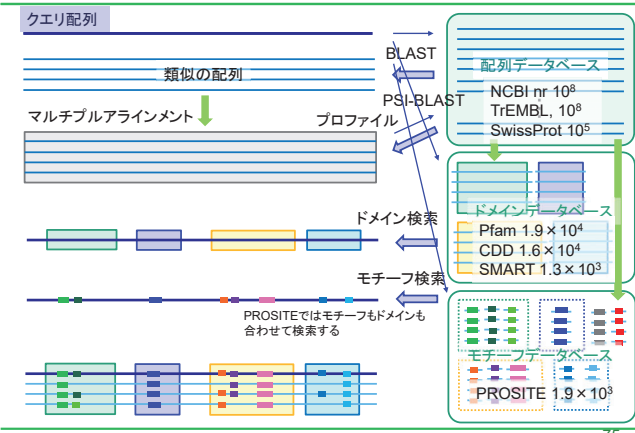
Name	Accession	Description	Start	End
PTKc_Src	cd05071	Catalytic domain of the Protein Tyrosine Kinase, Src. PTKs catalyze the transfer of the gamma phosphate of ATP to cytosolic (or non-receptor) PTKs, containing an SH2 domain, and is negatively regulated by phosphorylation at the C-terminal domain, forming of the oncogenic oncoprotein Src. Src is a member of the Src family of the oncogenic protein tyrosine kinases. The SH2 domain is a conserved domain of the Src family of the oncogenic protein tyrosine kinases. The SH2 domain is a conserved domain of the Src family of the oncogenic protein tyrosine kinases.	250-258	504-510

InterPro

- タンパク質の機能を解析するための総合的なサイト
- タンパク質のファミリー、ドメイン、重要な部位を示す配列特徴 (signature) を、複数のデータベースに対してまとめて検索する
- <https://www.ebi.ac.uk/interpro/>

The screenshot shows the InterPro website interface. The main heading is 'Classification of protein families'. Below the heading, there is a search bar and a 'Search by sequence' button. The search results show a list of protein families, including 'PTKc_Src'.

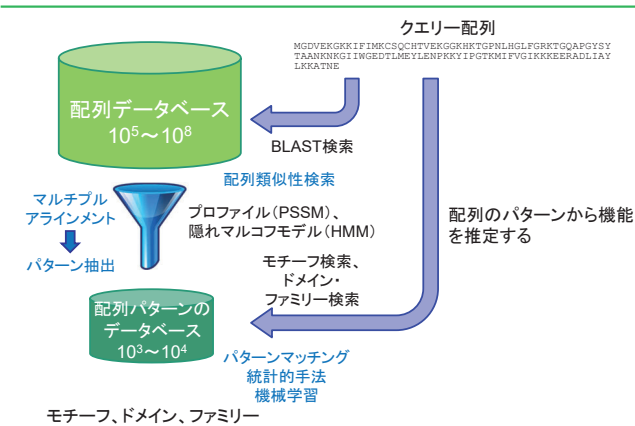
タンパク質の配列からの機能予測



InterProの検索(1)

The screenshot shows the InterPro search results for 'seq1.txt'. The search results show a list of protein families, including 'PTKc_Src'. A red arrow points to the 'FASTA形式の入力が要求され、配列のみだと、コメント行が自動的に追加される' (FASTA format input is required, only the sequence is provided, and comment lines are automatically added). Another red arrow points to the 'チロシンキナーゼSRCのアミノ酸配列 (seq1.txt) を入力して、「Search」を押す' (Enter the amino acid sequence of Tyrosine Kinase SRC (seq1.txt) and click 'Search').

配列から機能を予測する



InterProの検索(2)

The screenshot shows the InterPro search results for 'seq1.txt'. The search results show a list of protein families, including 'PTKc_Src'. A red arrow points to the '結果が得られるまで数分かかるかもしれない' (It may take a few minutes to get the results). Another red arrow points to the 'クリックして内容を見る' (Click to view content) button. A third red arrow points to the 'チェックが入ったら終了' (Check the box to finish) button.

InterProの検索(3)

Entry matches to this protein

複数のデータベースの検索結果がまとめて表示される

Domain: ドメイン構成

スーパーファミリーの情報

Active Site: 結合部位と活性部位

Binding Site: InterProに未統合のシグニチャ

モチーフ 81

UniProtKBの検索(1)

UniProtKB BLAST

How to use this tool

1. Enter either a protein or nucleotide sequence or a UniProt Identifier (e.g. G0D705 or A4_HUMAN or UP000000001) into the form field.

2. Optionally, change the program parameters with the dropdown menu under the form.

3. Click the Run BLAST button.

チロシナーゼSRCのアミノ酸配列(「seq1.txt」を入力して、「Search」を押す)

Run BLAST

モチーフ 85

InterProの検索(4)

Entry matches to this protein

SH3 domain

ドメイン名をクリックすると...

それぞれのバーにカーソルを当てると、ドメインの情報が表示される

ドメイン構成

モチーフ 82

UniProtKBの検索(2)

UniProtKB BLAST

検索条件が表示される

配列一致度に応じて色分け

「seq1.txt」の配列はSRC_HUMAN

検索結果のダウンロードとしてFASTA形式で検索された配列を列挙

モチーフ 86

参考

InterProの検索(5)

InterPro Classification of protein families

ドメイン構成の実例

タンパク質の実例

このドメインをもつタンパク質の例、既知の構造、バズウェイ、文献、他のデータベースへのリンクなど

このドメインに対応する他のデータベースのsignature

モチーフ 83

UniProtKBの検索(3)

UniProtKB - P12931 (SRC_HUMAN)

タンパク質名

遺伝子名

由来生物

機能

モチーフ 87

UniProtKB

- UniProt (Universal Protein resource): EBI, SBI, PIR (Protein Information Resources) が参加するコンソーシアム
- UniProtが提供するアノテーション付きのアミノ酸配列データベース
- <https://www.uniprot.org/>

UniProt

UniProtKB

UniRef

UniParc

Proteomes

モチーフ 84

UniProtKBの検索(4)

UniProtKB

触媒活性

さらし下の方を見ると...

モチーフ 88

UniProtKBの検索(5)

UniProtKBの検索(9)

UniProtKBの検索(6)

UniProtKBの検索(10)

UniProtKBの検索(7)

UniProtKBの検索(11)

UniProtKBの検索(8)

UniProtKBの検索(12)

UniProtKBの検索 (13)

UniProtKB Clustal Omegaによるアラインメントの結果が表示される

How to use this tool: align two or more protein sequences with the Clustal Omega program (see also this FAQ) to view their characteristics alongside each other.

1. Enter either protein sequences in FASTA format or UniProt identifiers into the form field, for example: TRN_HUMAN, TRN_PSS

2. Click the Run Align button.

Alignment Job status: COMPLETED

Download | Print and readout

How to print an alignment in color

Highlight

Annotations

- 1 Natural variant
- 2 Chain
- 3 Lipidation
- 4 Post-translational modification
- 5 Turn
- 6 Binding site
- 7 Inhibitor site
- 8 Signal peptide
- 9 Membrane
- 10 Transmembrane
- 11 Helical region
- 12 Disulfide bond
- 13 Hydrophobic region
- 14 Topological domain
- 15 Hydrophobic strand
- 16 Beta strand
- 17 Disordered region
- 18 Domain
- 19 Nucleotide binding
- 20 Active site
- 21 Similarity

モチーフ 97

ドメイン構成の比較

P12931のPROSITEの検索結果

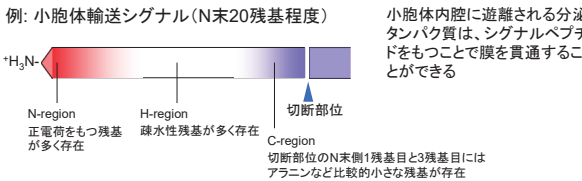
Q01973のPROSITEの検索結果

protein kinase domain

モチーフ 98

シグナル配列

- 細胞内局在、翻訳後修飾(リン酸化、糖鎖付加など)、転写因子結合など、特定の機能に関わる配列部位
- 3~60残基くらいのペプチド配列
- 必ずしも類縁でないタンパク質間で共通に見られる
- 厳密なコンセンサスは必ずしもないが、ある程度の配列特徴がある



典型的なシグナル配列の例

シグナル配列の機能	シグナル配列の例
核への輸送	-Pro-Pro-Lys-Lys-Lys-Arg-Lys-Val-
核からの輸送	-Leu-Ala-Leu-Lys-Leu-Ala-Gly-Leu-Asp-Ile-
ミトコンドリア内への輸送	*H ₃ N-Met-Leu-Ser-Leu-Arg-Gln-Ser-Ile-Arg-Phe-Phe-Lys-Pro-Ala-Thr-Arg-Thr-Leu-Cys-Ser-Ser-Arg-Tyr-Leu-Leu-
色素体内への輸送	*H ₃ N-Met-Val-Ala-Met-Ala-Met-Ala-Ser-Leu-Gln-Ser-Ser-Met-Ser-Ser-Leu-Ser-Ser-Ser-Asn-Ser-Phe-Leu-Gly-Gln-Pro-Leu-Ser-Pro-Ile-Thr-Leu-Ser-Pro-Phe-Leu-Gln-Gly-
ペルオキシソーム内への輸送	-Ser-Lys-Leu-COO'
小胞体内腔への輸送	*H ₃ N-Met-Met-Ser-Phe-Val-Ser-Leu-Leu-Leu-Val-Gly-Ile-Leu-Phe-Trp-Ala-Thr-Glu-Ala-Gln-Leu-Thr-Lys-Cys-Glu-Val-Phe-Gln-
小胞体内腔に保持	-Lys-Asp-Glu-Leu-COO'

重要な残基は色づけ → 赤: 正電荷、緑: 負電荷、灰色: 疎水性、青: OH基をもつ

モチーフ 100

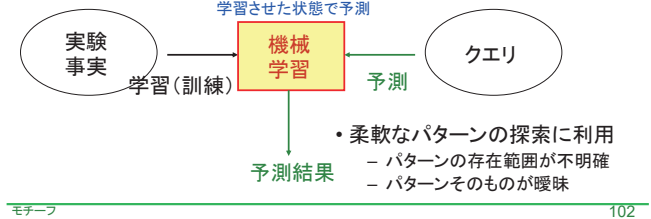
細胞内局在予測とシグナルペプチド予測

- 細胞内局在予測
 - TargetP
 - 真核生物のタンパク質の細胞内局在を予測
 - <http://www.cbs.dtu.dk/services/TargetP/>
 - PSORT
 - タンパク質の細胞内局在を予測
 - 生物種により異なるツールを用意
 - <https://psort.hgc.jp/>
- シグナルペプチド予測
 - SignalP
 - シグナルペプチドの存在および切断部位を予測
 - <http://www.cbs.dtu.dk/services/SignalP/>

モチーフ 101

機械学習の利用

- 機械学習 (machine learning)
 - 与えられたデータから知識を習得し、判別や分類を行う
 - 主な手法
 - ニューラルネットワーク
 - HMM (Hidden Markov Model)
 - SVM (Support Vector Machine)
 - ランダムフォレスト
 - 深層学習 など
 - 応用
 - モチーフ検索、二次構造予測、翻訳後修飾予測、相互作用部位予測



TargetPの利用

TargetP-2.0 Server <http://www.cbs.dtu.dk/services/TargetP/>

TargetP-2.0: Click here to run TargetP-2.0

Predict: Instructions/help | Data | Applications | Protein version

Submit data

TargetP-2.0 server predicts the presence of N-terminal sequences: signal peptide (SP), mitochondrial transit peptide (mTP), chloroplast transit peptide (cTP) or thylakoid lumenal transit peptide (lTP). For the sequences predicted to contain an N-terminal sequence a potential cleavage site is also predicted.

Paste or upload protein sequence(s) as fasta format.

Protein sequences should be not less than 10 amino acids. The maximum number of proteins is 5000.

seq4.txtの配列を入力

「Plant」を選択

「Submit」ボタンを押す

モチーフ 103

TargetPの予測結果の例

TargetP-2.0

Summary of 1 predicted sequences from Plant

Predictions list. Use the help page for more detailed description of the output page.

Predicted proteins

葉緑体輸送ペプチド (スコア最大)

ミトコンドリア輸送ペプチド

局在場所の判定

チラコイド内腔輸送ペプチド

cs: cleavage site

Protein type	Signal peptide	Mitochondrial transfer peptide	Chloroplast transfer peptide	Thylakoid lumenal transfer peptide
Likelihood	0.001	0	0.9974	0.0014

モチーフ 104

課題(2/2)

- 講義のページのファイル`seqx.txt`のヒトのタンパク質の配列について、以下の問いに答えよ。
 - 適当な予測システムを用いて、このタンパク質がシグナルペプチドをもつかもたないかを予測せよ。もつと予測される場合、その位置(残基番号)を求めよ。使用した予測システムを記すこと。
 - 適当な予測システムを用いて、このタンパク質が膜貫通領域をもつかもたないかを予測せよ。もつと予測される場合、その位置(残基番号)を求めよ。使用した予測システムを記すこと。
 - このタンパク質が登録されている配列データベース([UniProtKB](#)など)のアノテーション情報と、(1)、(2)の予測結果を比較せよ。

課題

113

Attendance check and assignments submission

- Attendance is taken via ITC-LMS.
 - You can find one-time password in the zoom chat.
- Submit the assignments to the ITC-LMS site.
 - If you cannot use ITC-LMS, you can submit the assignments to the following address via e-mail.
 - shimizu@bi.a.u-tokyo.ac.jp
- The class ends at 20:30. You should submit the assignments before the end of the class.
 - If you cannot meet the deadline, please contact us via the zoom chat.

課題

117

出席の確認と課題の提出

- 出席は、ITC-LMSから行って下さい。
 - ワンタイムパスワードは、授業開始時にzoomのチャットでお伝えします。
- 課題の提出は、ITC-LMSから行って下さい。
 - ITC-LMSから提出できない場合は以下のメールアドレスに送って下さい。
 - shimizu@bi.a.u-tokyo.ac.jp
- 授業は20:30に終了します。課題は、そのときまで提出して下さい。
 - 間に合わない場合はzoomのチャットで連絡して下さい。
 - 4月23日まで(24日0時まで)に提出して下さい。
 - Submit the assignments until April 23.

課題

114

Assignments (1 / 2)

- The purine repressor protein, PurR, is a member of the lac repressor family of Escherichia coli DNA-binding proteins (UniProtKB ID: [P0ACP7 \(PURR ECOLI\)](#)). Answer the following questions.
 - Describe the name, the pattern (regular expression), and the position (in residue number) of the motif which is related to DNA binding contained in PurR. You can use [PROSITE](#).
 - This protein is known to have sequences similar to those of E. coli periplasmic D-galactose binding proteins (DBPs). Perform a PSI-BLAST search against Swiss-Prot database in the [NCBI BLAST site](#) and describe the UniProtKB ID of the DBP. Also describe the aligned region of the two proteins and the sequence identity of the region.
 - Answer whether the above motif (1) is aligned to the DBP of (2) or not when using PSI-BLAST.

Swiss-Prot sequences are identified by the UniProtKB IDs.

課題

115

Assignments (2 / 2)

- For the human protein sequence in "`seqx.txt`", answer the following questions.
 - Is this protein predicted to contain a signal peptide? If so, describe its position (in residue number). Use an appropriate prediction tool.
 - Is this protein predicted to contain transmembrane regions? If so, describe their position (in residue number). Use an appropriate prediction tool.
 - Compare the annotation described in the sequence database (e.g. [UniProtKB](#)) with the prediction results of (1) and (2).

課題

116